

US009319352B1

(12) **United States Patent**
Nielsen

(10) **Patent No.:** **US 9,319,352 B1**
(45) **Date of Patent:** ***Apr. 19, 2016**

(54) **EFFICIENT MESSAGE SWITCHING IN A SWITCHING APPARATUS**

(71) Applicant: **Marvell International Ltd.**, Hamilton (BM)

(72) Inventor: **Jacob V. Nielsen**, Virum (DK)

(73) Assignee: **Marvell International Ltd.**, Hamilton (BM)

(*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 243 days.
This patent is subject to a terminal disclaimer.

4,536,873 A	8/1985	Leete
4,703,475 A	10/1987	Dretzka et al.
4,787,081 A	11/1988	Waters et al.
4,807,408 A	2/1989	Lew et al.
4,816,826 A	3/1989	Munter et al.
4,821,259 A	4/1989	DeBruler et al.
4,831,373 A	5/1989	Hess
4,837,858 A	6/1989	Ablay et al.
4,953,158 A	8/1990	Schrur
5,058,110 A	10/1991	Beach et al.
5,093,827 A	3/1992	Franklin et al.

(Continued)

FOREIGN PATENT DOCUMENTS

EP	1 206 087 A2	5/2002
WO	WO-2004049174 A2	6/2004

OTHER PUBLICATIONS

(21) Appl. No.: **14/053,128**

(22) Filed: **Oct. 14, 2013**

U.S. Appl. No. 11/187,236, filed Jul. 22, 2005, Nielsen, et al.
(Continued)

Related U.S. Application Data

(63) Continuation of application No. 11/187,236, filed on Jul. 22, 2005, now Pat. No. 8,559,443.

(51) **Int. Cl.**
H04L 12/58 (2006.01)
H04L 12/931 (2013.01)

(52) **U.S. Cl.**
CPC **H04L 49/60** (2013.01)

(58) **Field of Classification Search**
USPC 370/424, 331, 347
See application file for complete search history.

(56) **References Cited**

U.S. PATENT DOCUMENTS

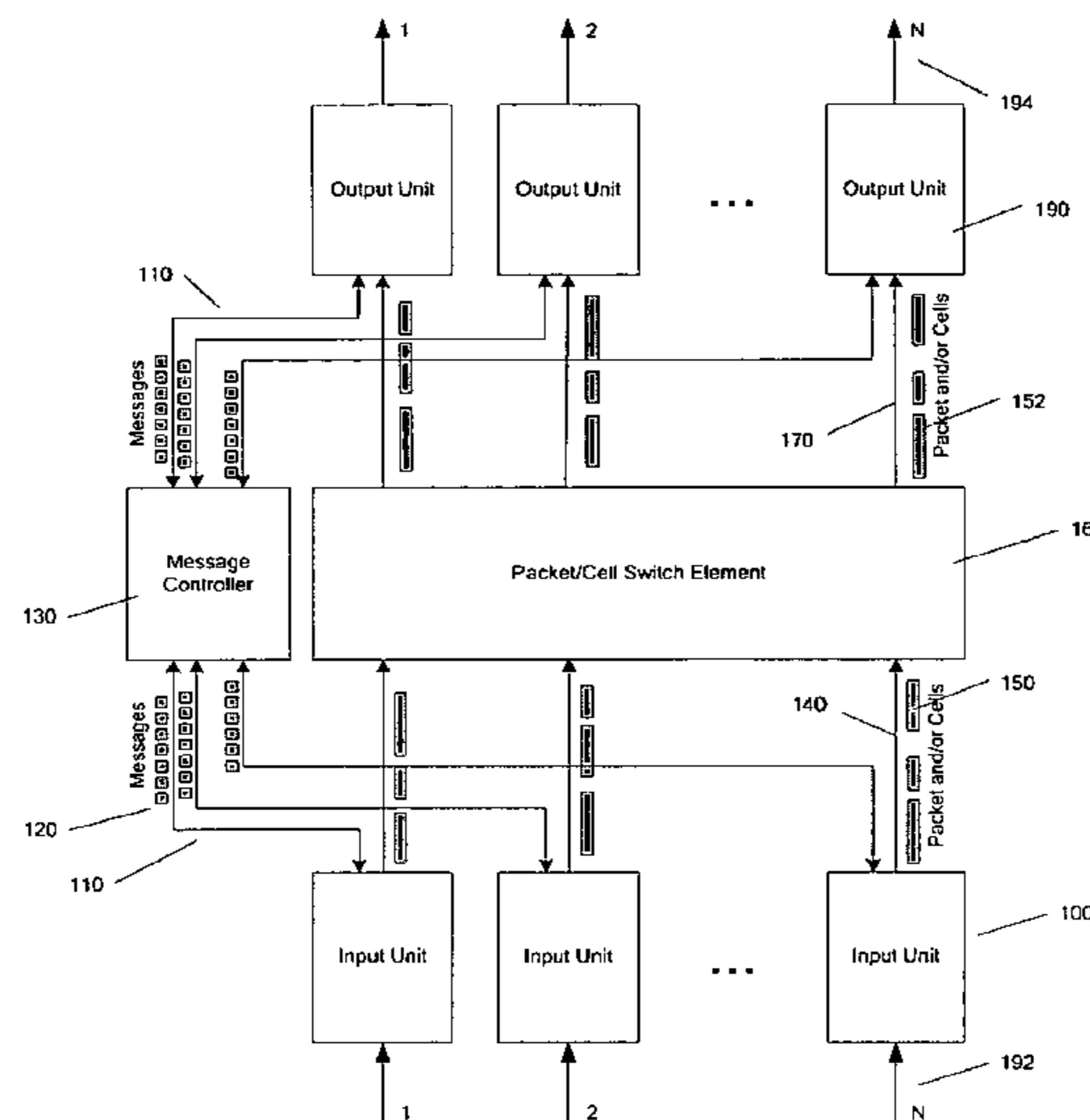
4,156,798 A	5/1979	Doelz
4,271,506 A	6/1981	Broc et al.
4,387,425 A	6/1983	El-Gohary
4,475,010 A	10/1984	Huensch et al.

Primary Examiner — Dang Ton
Assistant Examiner — Pamit Kaur

(57) **ABSTRACT**

A switching apparatus including a switch element, a message controller separate from the switch element, a plurality of input units, and a plurality of output units. The plurality of input units are configured to provide packets to the switch element and provide, to the message controller, messages corresponding to respective ones of the packets. The plurality of output units are configured to receive, from the message controller, the messages, and receive, from the switch element, the packets to be output from the switching apparatus. The switch element is configured to transfer each of the packets between the input units and the output units. The message controller is configured to transfer each of the messages corresponding to the respective ones of the packets between the input units and the output units.

16 Claims, 6 Drawing Sheets



(56)

References Cited

U.S. PATENT DOCUMENTS			
5,129,077	A	7/1992	Hillis
5,184,347	A	2/1993	Farwell et al.
5,218,367	A	6/1993	Sheffer et al.
5,227,775	A	7/1993	Bruckert et al.
5,239,537	A	8/1993	Sakauchi
5,239,680	A	8/1993	Grube et al.
5,251,212	A	10/1993	Gass
5,282,201	A	1/1994	Frank et al.
5,296,934	A	3/1994	Ohtsuki
5,341,483	A	8/1994	Frank et al.
5,355,370	A	10/1994	Fairhurst et al.
5,367,520	A	11/1994	Cordell
5,387,905	A	2/1995	Grube et al.
5,388,243	A	2/1995	Glider et al.
5,396,359	A	3/1995	Abramovitz
5,404,374	A	4/1995	Mullins et al.
5,404,537	A	4/1995	Olnowich et al.
5,408,419	A	4/1995	Wong
5,420,574	A	5/1995	Erickson et al.
5,422,816	A	6/1995	Sprague et al.
5,425,025	A	6/1995	Tahara
5,432,908	A	7/1995	Heddes et al.
5,434,993	A	7/1995	Liencrest et al.
5,442,659	A	8/1995	Bauchot et al.
5,481,537	A	1/1996	Crisler et al.
5,491,692	A	2/1996	Gunner et al.
5,493,569	A	2/1996	Buchholz et al.
5,500,858	A	3/1996	McKeown
5,515,379	A	5/1996	Crisler et al.
5,519,706	A	5/1996	Bantz et al.
5,530,693	A	6/1996	Averbuch et al.
5,548,814	A	8/1996	Lorang et al.
5,548,818	A	8/1996	Sawyer et al.
5,594,868	A	1/1997	Nakagoshi et al.
5,598,417	A	1/1997	Crisler et al.
5,600,643	A	2/1997	Robrock, II
5,666,364	A	9/1997	Pierce et al.
5,668,803	A	9/1997	Tymes et al.
5,675,807	A	10/1997	Iswandhi et al.
5,678,172	A	10/1997	Dinkins
5,696,760	A	12/1997	Hardin et al.
5,703,875	A	12/1997	Burnett
5,737,706	A	4/1998	Seazholtz et al.
5,748,627	A	5/1998	Weir
5,751,955	A	5/1998	Sonnier et al.
5,761,193	A	6/1998	Derango et al.
5,787,075	A	7/1998	Uchida
5,790,963	A	8/1998	Welham
5,799,012	A	8/1998	Ayerst et al.
5,809,012	A	9/1998	Takase et al.
5,809,076	A	9/1998	Hofmann
5,835,485	A	11/1998	Grube et al.
5,838,894	A	11/1998	Horst
5,872,523	A	2/1999	Dellaverson et al.
5,883,893	A	3/1999	Rumer et al.
5,901,142	A	5/1999	Averbuch et al.
5,913,028	A	6/1999	Wang et al.
5,914,953	A	6/1999	Krause et al.
5,931,901	A	8/1999	Wolfe et al.
5,963,542	A	10/1999	Shum
5,963,555	A	10/1999	Takase et al.
5,974,458	A	10/1999	Abe et al.
5,987,018	A	11/1999	Freeburg et al.
6,002,678	A	12/1999	Jayapala et al.
6,028,860	A	2/2000	Laubach et al.
6,038,591	A	3/2000	Wolfe et al.
6,041,038	A	3/2000	Aimoto
6,044,061	A	3/2000	Aybay et al.
6,058,102	A	5/2000	Drysdale et al.
6,067,297	A	5/2000	Beach
6,069,886	A	5/2000	Ayerst et al.
6,075,784	A	6/2000	Frankel et al.
6,081,512	A	6/2000	Muller et al.
6,085,233	A	7/2000	Jeffrey et al.
6,137,798	A	10/2000	Nishihara et al.
6,147,980	A	11/2000	Yee et al.
6,157,843	A	12/2000	Derango et al.
6,161,142	A	12/2000	Wolfe et al.
6,167,457	A	12/2000	Eidison et al.
6,185,213	B1	2/2001	Katsube et al.
6,188,892	B1	2/2001	Krishnamurthi et al.
6,198,927	B1	3/2001	Wright et al.
6,198,929	B1	3/2001	Krishnamurthi et al.
6,240,067	B1	5/2001	Sorber
6,240,444	B1	5/2001	Fin et al.
6,246,691	B1	6/2001	Briem et al.
6,249,528	B1	6/2001	Kothary
6,272,141	B1	8/2001	Reed
6,272,338	B1	8/2001	Modzelesky et al.
6,272,341	B1	8/2001	Threadgill et al.
6,272,581	B1	8/2001	Leung et al.
6,278,861	B1	8/2001	Ward et al.
6,333,932	B1	12/2001	Kobayasi et al.
6,351,466	B1	2/2002	Prabhakar et al.
6,366,622	B1	4/2002	Brown et al.
6,370,127	B1	4/2002	Daraiseh et al.
6,370,153	B1	4/2002	Eng
6,370,381	B1	4/2002	Minnick et al.
6,373,399	B1	4/2002	Johnson et al.
6,385,174	B1	5/2002	Li
6,389,022	B1	5/2002	Jeong et al.
6,389,479	B1	5/2002	Boucher et al.
6,392,994	B1	5/2002	Dubuc
6,402,691	B1	6/2002	Peddicord et al.
6,411,620	B1	6/2002	Takase et al.
6,421,357	B1	7/2002	Hall
6,452,926	B1	9/2002	Wiklund
6,512,741	B1	1/2003	Kohzuki et al.
	H	2/2003	Ledsham et al.
6,515,991	B1	2/2003	McKeown
6,529,848	B2	3/2003	Sone
6,545,996	B1	4/2003	Falco et al.
6,574,211	B2	6/2003	Padovani et al.
6,621,796	B1	9/2003	Miklos
6,628,615	B1	9/2003	Joseph et al.
6,631,180	B1	10/2003	Ikeda
6,639,916	B1	10/2003	Wakizaka
6,697,349	B2	2/2004	Mathis et al.
6,704,394	B1	3/2004	Kambhatla et al.
6,731,638	B1	5/2004	Ofek
6,738,381	B1	5/2004	Agnevik et al.
6,751,232	B1	6/2004	Patterson et al.
6,778,512	B2	8/2004	Gipson et al.
6,907,001	B1	6/2005	Nakayama et al.
6,907,254	B1	6/2005	Westfield
6,928,304	B2	8/2005	Wigell et al.
6,944,138	B1	9/2005	Song
6,954,448	B2	10/2005	Farely et al.
6,956,818	B1	10/2005	Thodiyil
6,970,478	B1	11/2005	Nishihara
6,977,919	B1	12/2005	Stanwood
7,065,580	B1	6/2006	Eberie et al.
7,079,550	B2	7/2006	Padovani et al.
7,088,710	B1	8/2006	Johnson et al.
7,154,854	B1	12/2006	Zweig et al.
7,173,922	B2	2/2007	Beach
7,177,309	B2	2/2007	Shinohara
7,203,193	B2	4/2007	Hoof
7,277,425	B1	10/2007	Sikdar
7,280,511	B2	10/2007	Ahn
7,362,751	B2	4/2008	Khacherian et al.
7,742,486	B2	6/2010	Nielsen et al.
8,559,443	B2*	10/2013	Nielsen 370/395.4
2001/0021174	A1	9/2001	Luijten et al.
2001/0055283	A1	12/2001	Beach
2002/0051427	A1	5/2002	Carvey
2002/0054568	A1	5/2002	Hoogenboom et al.
2002/0080775	A1	6/2002	Engbersen et al.
2002/0101839	A1	8/2002	Farley et al.
2002/0118689	A1	8/2002	Luijten et al.
2002/0145974	A1	10/2002	Saidi et al.
2002/0152263	A1	10/2002	Goldrian et al.
2002/0163915	A1	11/2002	Wallner et al.
2002/0170013	A1	11/2002	Bolourchi et al.

(56)

References Cited

U.S. PATENT DOCUMENTS

2002/0181455 A1 12/2002 Norman et al.
2003/0002865 A1 1/2003 Matsui et al.
2003/0016688 A1 1/2003 Hoof
2003/0021230 A1 1/2003 Kuo et al.
2003/0043772 A1 3/2003 Mathis et al.
2003/0056073 A1 3/2003 Zeiger
2003/0058802 A1 3/2003 Jones et al.
2003/0060203 A1 3/2003 Ahn
2003/0063583 A1 4/2003 Padovani et al.
2003/0063604 A1 4/2003 Wallner et al.
2003/0063618 A1 4/2003 Khacherian et al.
2003/0088694 A1 5/2003 Patek et al.
2003/0103501 A1 6/2003 Clem et al.
2003/0112820 A1 6/2003 Beach
2003/0117947 A1 6/2003 Koehl et al.
2003/0117983 A1 6/2003 Ton et al.
2003/0118058 A1 6/2003 Kim et al.
2003/0123469 A1 7/2003 Nong

2003/0221043 A1 11/2003 Sota
2003/0227932 A1 12/2003 Meempat et al.
2003/0233612 A1 12/2003 Gilchrist et al.
2004/0085979 A1 5/2004 Lee et al.
2004/0093415 A1 5/2004 Thomas
2004/0148396 A1 7/2004 Meyer et al.
2004/0268057 A1 12/2004 Landin et al.
2005/0008011 A1 1/2005 Georgiou et al.
2005/0041970 A1 2/2005 Harai
2005/0047334 A1 3/2005 Paul et al.
2005/0132060 A1 6/2005 Mo et al.
2005/0135356 A1 6/2005 Muthukrishnan et al.
2005/0135398 A1 6/2005 Muthukrishnan et al.
2006/0031506 A1 2/2006 Redgate
2006/0285520 A1 12/2006 Venkitaraman
2007/0019674 A1 1/2007 Bourlas et al.

OTHER PUBLICATIONS

U.S. Appl. No. 10/898,540, filed Jul. 2, 2004, Nielson, et al.

* cited by examiner

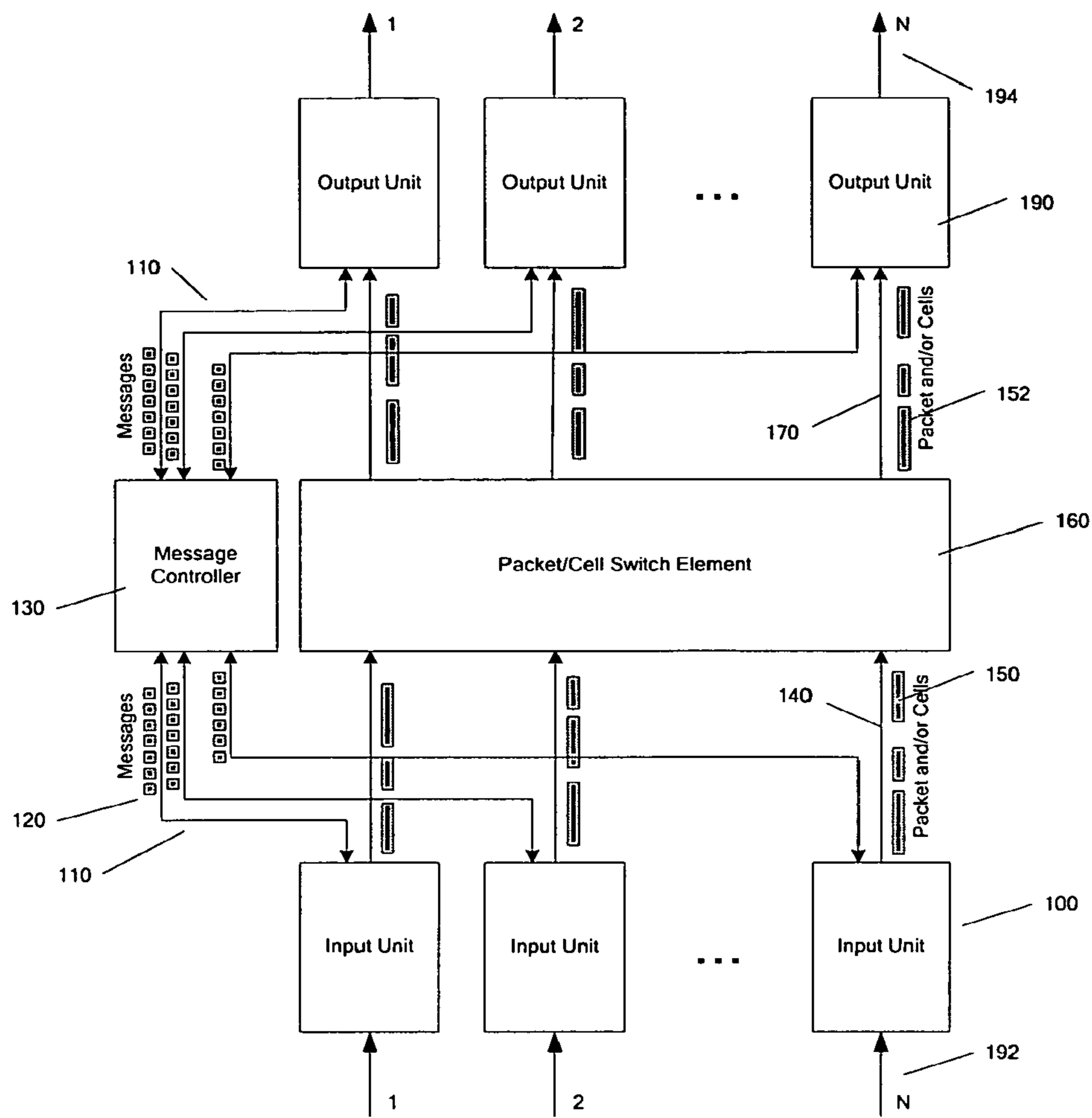


FIGURE 1

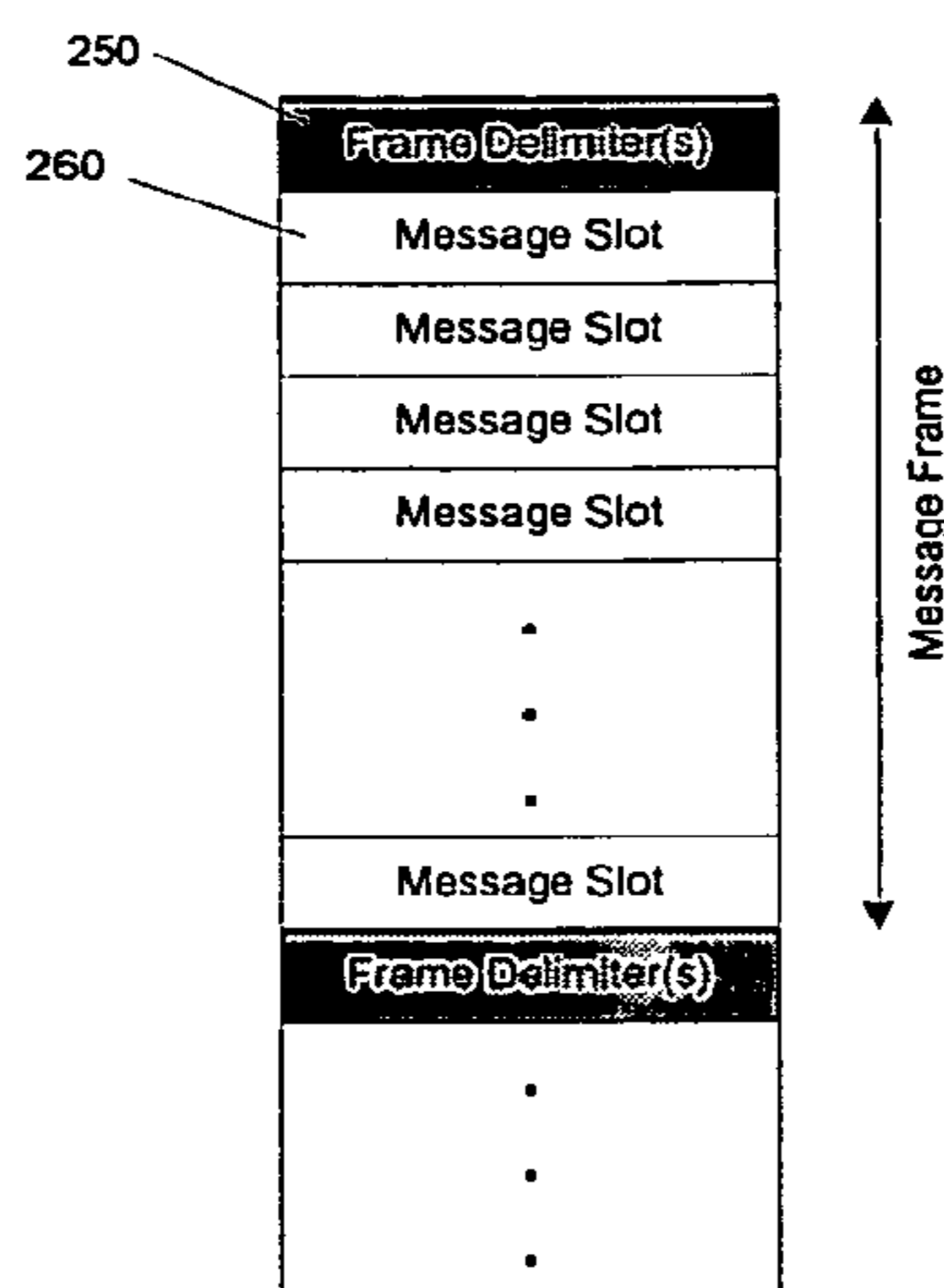
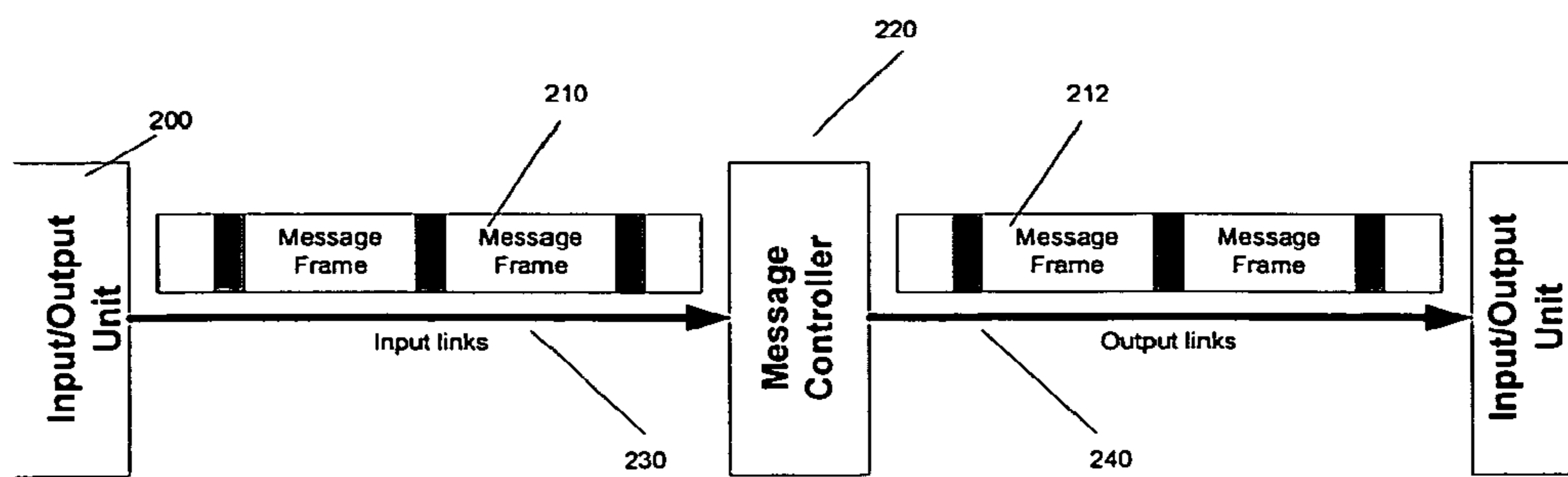


FIGURE 2

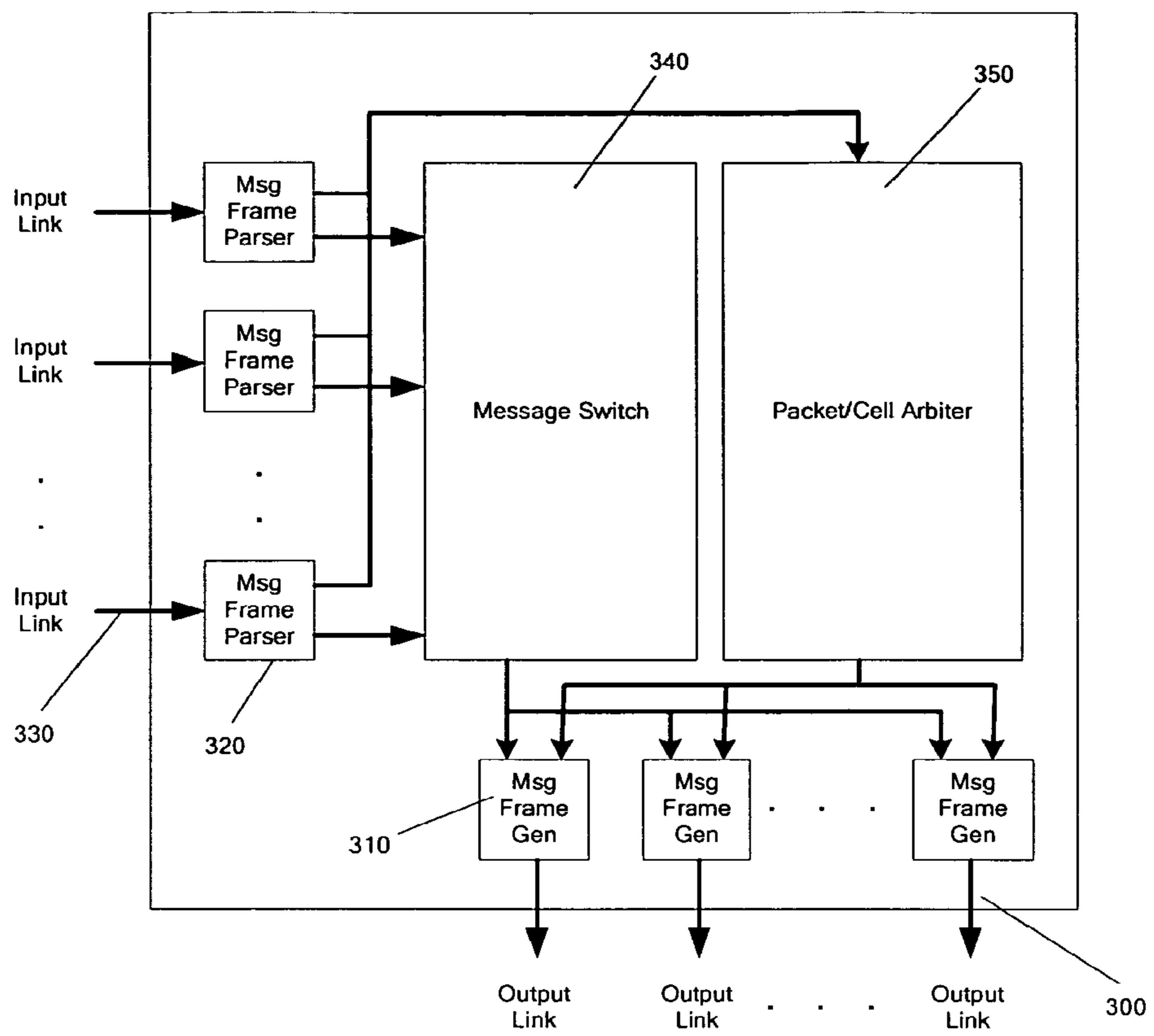


FIGURE 3

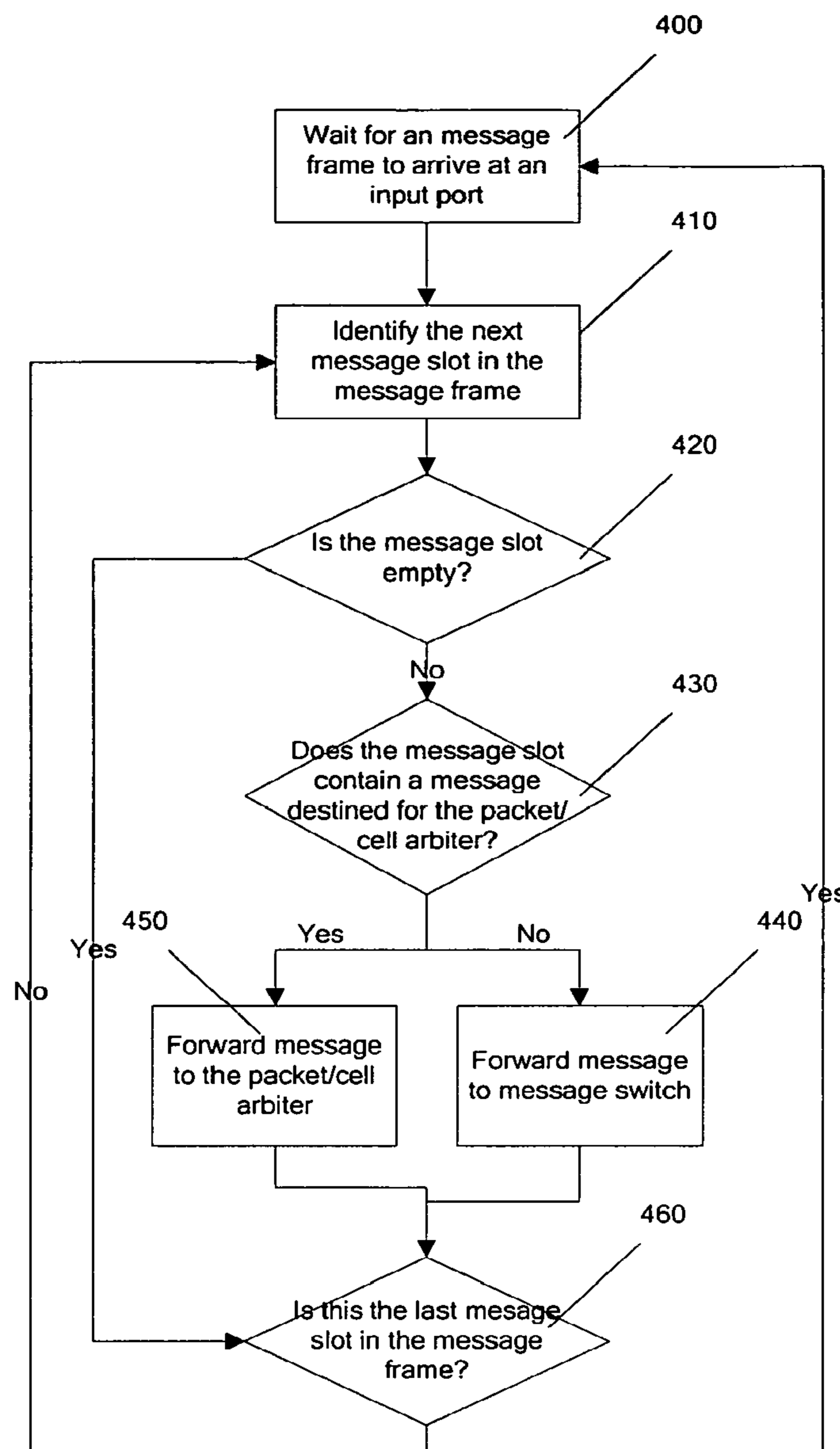


FIGURE 4

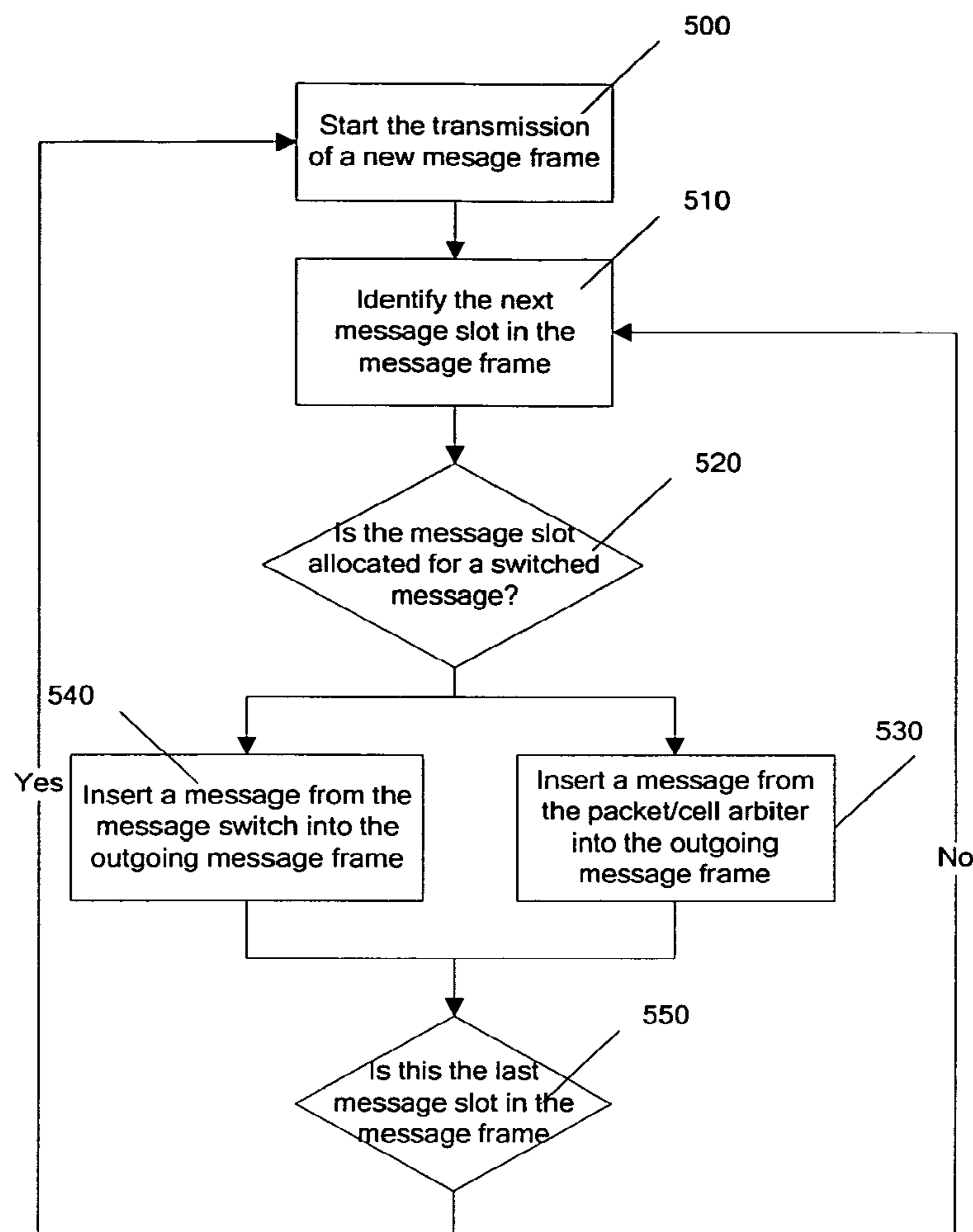


FIGURE 5

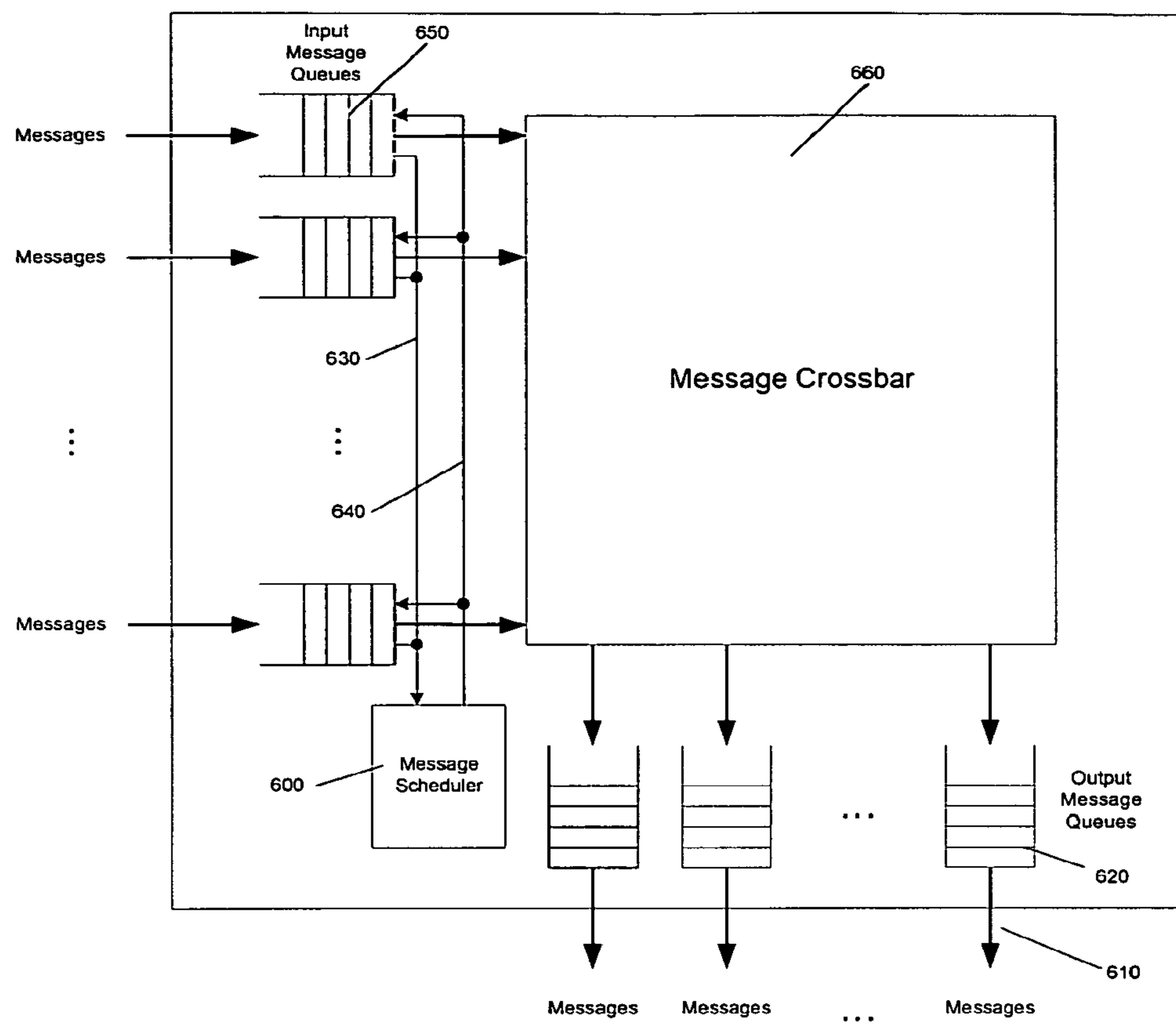


FIGURE 6

EFFICIENT MESSAGE SWITCHING IN A SWITCHING APPARATUS

CROSS-REFERENCE TO RELATED APPLICATIONS

The present disclosure is a continuation of U.S. patent application Ser. No. 11/187,236 filed on Jul. 22, 2005. The entire disclosure of the application referenced above is incorporated herein by reference.

FIELD

The invention relates generally to switching of messages in a packet/cell switching apparatus. The message switching is optimized for efficient switching of small messages, and is performed fully separated from the packet/cell switching.

BACKGROUND

The background description provided herein is for the purpose of generally presenting the context of the disclosure. Work of the presently named inventors, to the extent the work is described in this background section, as well as aspects of the description that may not otherwise qualify as prior art at the time of filing, are neither expressly nor impliedly admitted as prior art against the present disclosure.

The input and output units of a general switching apparatus are connected respectively to the input and output links of a packet/cell switch element which resides inside the switching apparatus. Incoming packets/cells are switched from the input units to the output units via the packet/cell switch element in a packet and/or cell format. The typical packet format is a variable size frame with a typical size range from 32 to 10000 bytes, and the typical cell format is a fixed size frame with a typical size range from 32 to 80 bytes.

The input and output units of the switching apparatus may also require means for efficiently switching messages between the input and output units. Such messages are typically used for distributing information related to packet/cell input/output unit queuing status, packet/cell switching/scheduling credits, packet/cell flow control commands, and packet/cell control table state information. These messages are typically in the range of 2 to 16 bytes, which is smaller than the typical minimum packets and/or cell size. Furthermore, messages may be switched from input units to output units and vice versa, while packets/cells are typically only switched from input units to output units.

The packet/cell switch element is typically optimized for switching of packet and/or cells with a minimum size of 32 to 80 bytes, and therefore is inefficient for switching the smaller messages. One reason for this inefficiency is that a required switching header per packet/cell unit may be comparable in size to the message itself. The packet/cell switch element may also pad the size of the message up to a minimum packet/cell size, which also reduces the efficiency of the packet/cell switch element when used to switch small messages. It may also be a problem that when messages are switched across the packet/cell switch element together with packets/cells, the messages impact the packet/cell switching throughput and vice versa, and this results in non-deterministic switching performance for both messages and packets/cells.

One solution described in U.S. Patent Publication No. 2003/0103501 uses a separate ring element integrated inside a switch element to separate smaller messages from traffic data (packets/cells) which is switched across a crossbar. The ring element is constructed by successively connecting adja-

cent switch element links, forming a ring for passing the messages from an input link, successively through intermediate links, to the destination output link. The drawback of this approach is that although the messages and traffic data (packets/cells) use separate switching resources inside the switch element, they share the switch elements input and output links when passed to and from the switch element respectively. This structure means that the messages impact the switching of traffic data (packets/cells) and vice versa, which may result in non-deterministic switching performance for messages and traffic data (packets/cells).

Another solution described in U.S. Pat. No. 5,703,875 uses separate queuing resources inside a switch element to separate short control messages from longer data messages. Each input link has separate input queue resources to separate short and long messages, and all messages are switched using the same crossbar element. The drawback of this approach is that although the messages and traffic data (packets/cells) use separate queue resources inside the switch element, they share the switch elements input and output links when passed to and from the switch element respectively, and they also share the same crossbar element. This structure means that, the messages impact the switching of traffic data and vice versa, which may result in non-deterministic switching performance for messages and traffic data (packets/cells).

SUMMARY

At least one aspect of the present invention performs efficient message switching inside a packet/cell switching apparatus, fully separated from the packet/cell switching.

According to one aspect of the invention, there is provided a method of transferring packets/cells and messages within a switching apparatus that includes a plurality of input units, a packet/cell switch element, a message controller, and a plurality of output units. The method includes generating a message at one of the plurality of input units and output units, the message destined for another of the input units and output units. The method also includes transferring the message, via the message controller and via one of a plurality of links dedicated for message transfer, from the one of the plurality of input units and outputs units, to another of the input units and the output units. The method further includes outputting a packet/cell scheduling request command from the one of the input units to the message controller, the packet/cell scheduling request command being transferred to the message controller from one of the plurality of input units using one of a plurality of links dedicated for message transfer within the switching apparatus. The method still further includes receiving the packet/cell transfer scheduling request command at the message controller, determining by the message controller when to allow transfer of the packet/cell, and notifying the one of the plurality of input units by outputting a packet/cell data acknowledging command from the message controller to the one of the plurality of input units over the one or another of the plurality of links dedicated for message transfer. The method also includes outputting the packet/cell from the one of the plurality of input units to the packet/cell switch element, by using one of a plurality of links dedicated for packet/cell transfer.

According to another aspect of the invention, there is provided a method of transferring packets/cells and messages within a switching apparatus that includes a plurality of input units, a packet/cell switch element, a message controller, and a plurality of output units. The method includes generating a message at one of the plurality of input units and output units, the message destined for another of the input units and output

3

units. The method also includes transferring the message, via a message switch of the message controller and via at least one of a plurality of links dedicated for message transfer, from the one of the plurality of input units and outputs units, to another of the input units and the output units. The method further includes outputting a packet/cell from one of the input units to one of the output units, via the packet/cell switch element and via at least one of a plurality of links dedicated for packet/cell transfer, under control of the message controller.

According to yet another aspect of the invention, there is provided a system for transferring packets/cells and messages within a switching apparatus that includes a plurality of input units, a packet/cell switch element, a message controller which includes a packet/cell arbiter and a message switch, and a plurality of output units. The system includes a first plurality of input and output links for respectively connecting each of the input units and the output units to a packet/cell switch element. The system also includes a second plurality of input and output links for connecting each of the input and output units to a message controller. All packets and cells are transferred from the input units to the output units by way of the first plurality of input and output links and the packet/cell switch element, under packet/cell scheduling control including a first transfer of a scheduling request messages from one of the input units to the packet/cell arbiter by way of the second plurality of input links and then one of input units receiving a corresponding scheduling acknowledge messages back from the packet/cell arbiter by way of the second plurality of output links, and then a second transfer of the corresponding packet/cell to one of the output units by way of a first plurality of input and output links and the packet/cell switch element. All messages are transferred among the input units and the output units by way of the second plurality of input and output links, and the message switch.

Further areas of applicability of the present disclosure will become apparent from the detailed description, the claims and the drawings. The detailed description and specific examples are intended for purposes of illustration only and are not intended to limit the scope of the disclosure.

BRIEF DESCRIPTION OF DRAWINGS

The foregoing advantages and features of the invention will become apparent upon reference to the following detailed description and the accompanying drawings, of which:

FIG. 1 shows a switching apparatus according to a first embodiment of the invention.

FIG. 2 shows the transmission format of messages to and from the message controller according to the first embodiment.

FIG. 3 shows a top-level block diagram of the message controller which includes a packet/cell arbiter and a message switch according to the first embodiment.

FIG. 4 shows a flow chart for the message controller's message frame parser function according to the first embodiment.

FIG. 5 shows a flow chart for the message controller's message frame generator function according to the first embodiment.

FIG. 6 shows a block diagram of a message switch inside the message controller according to the first embodiment.

DESCRIPTION

A switching apparatus according to a first embodiment of the invention includes a packet/cell switch element and a

4

message controller. The switching apparatus enables efficient message switching via the message controller, fully separated from the packet/cell switching which is performed via the packet/cell switch element.

In addition to operating as a mechanism for enabling efficient message switching between the input and output units of the switching apparatus, the message controller performs a packet/cell scheduling arbitration by processing received requests messages from input units and generating and transmitting acknowledge messages back to input units for directing packets/cells across the packet/cell switch element.

The messages are transmitted to the message controller in a frame format. The frame format defines multiple message transmission timeslots per frame, and the position of each message transmission timeslot is fixed relative to the frame boundary. Although transmission delineation overhead is required per frame, it is not required per individual message transmission timeslot, thereby providing an efficient message transmission format with little overhead.

When a message arrives at the message controller, it will either be forwarded to the message controller's packet/cell arbiter or to the message controller's message switch, depending upon if the message is a packet/cell scheduling request message type or an message type which is to be exchanged between the input and output units of the switching apparatus respectively.

In a typical switching apparatus embodiment, the message controller's packet/cell arbiter will accept packet/cell scheduling request command messages and generate packet/cell scheduling acknowledge command message in return. The operation of the packet/cell arbiter is outside the scope of this invention, and will not be discussed in any detail herein, whereby the embodiments of the present invention are independent of the packet/cell arbiter.

The message controller integrates a message switch which is optimized for small messages. The message switch is typically optimized for smaller messages in the typical size range from 2-16 bytes. Since the message switch can be optimized for switching of very small messages independent of the packet/cell switch element, it is possible to integrate a highly efficient message switch.

The message controller's message switch includes a set of input message queues per input link, and a set of output message queues per output link, whereby these queues are connected via the inputs and outputs of a message crossbar, respectively. A message scheduler controls the switching of messages across the message crossbar, whereby the message crossbar is capable of simultaneously switching multiple messages from one or more input message queues, to one or more output message queues, on a per output link basis.

FIG. 1 shows components of an $N \times N$ switching apparatus according to a first embodiment of the invention (N is an integer value greater than one). The size of the switching apparatus may be other than $N \times N$, such as $M \times N$ where M and N are different integers. The $N \times N$ switching apparatus includes N input ports **192**, N output ports **194**, N input units **100**, N output units **190**, a message controller **130**, a packet/cell switch element **160**, message input/output links **110**, packet/cell input links **140**, and packet/cell output links **170**.

Each of the N input ports **192** receives packet and/or cells, and buffers them in their respective input units **100**, in a manner known to those skilled in the art. Each input unit **100** connects to a packet/cell switch element **160** via one or more input links **140**, and the packets and/or cells **150** are transmitted from the input units **100** to the packet/cell switch element **160** via these dedicated packet/cell input links **140**. Each output unit **190** connects to a packet/cell switch element **160**

5

via one or more output links **170**, to transmit packets and/or cells **152** from the packet/cell switch element **160** to the output units **190** via these dedicated packet/cell output links **170**, before final forwarding to their destination output port **194**.

In addition to the connectivity between the input/output units and the packet/cell switch element, each input unit **100** connects to the message controller **130** via one or more input/output links **110** that are dedicated for bi-directional transfer of messages between the message controller **130** and the input units **100** within the N×N switching apparatus. An input unit **100** transmits messages **120** to the message controller **130** via one or more input/output links **110** that is dedicated for message transfer, and receives messages **120** from the message controller **130** via one or more input/output links **110** that is dedicated for message transfer. Similarly, each output unit **190** also connects to the message controller **130** via one or more input/output links **110** that are dedicated for bi-directional message transfer between the message controller **130** and the output units **190**.

While FIG. **1** shows the input/output links **110** as bidirectional links (arrows at both ends), the individual links are preferably uni-directional with some of the links **110** dedicated for transfer of messages from the input/output units **100/190** to the message controller **130**, while other links are dedicated for transfer of messages from the message controller **130** to the input/output units **100/190**.

A packet/cell arbitration (scheduling) function is included in the message controller **130** for the embodiment structure shown in FIG. **1**. The packet/cell arbitration function processes received requests messages from input units and generates and transmits acknowledge messages to input units for directing packets/cells from the input units across the packet/cell switch element for switching to an output unit, before final forwarding to their destination output port. The packet/cell arbitration function included in the message controller may cooperate with packet/cell arbitration functions embodied in one or more of the input units **100**, one or more of the output units **190**, or any combination of these components, depending on the specific switching apparatus embodiment.

A preferred implementation of the packet/cell switch element **160** is a single stage structure of parallel switch devices, scheduled such that packets/cells from the input units are distributed in parallel across these parallel switch devices.

In a preferred embodiment of the first embodiment, input unit **L** and output unit **L** are integrated into a single physical device. This way, the integrated input and output unit can share the same input link connecting to the message controller **130**, which reduces the number of input/output links **110** on the message controller **130** by a factor of two.

FIG. **2** shows the transmission format for the input links **230** and output links **240** of the message controller **130** according to the first embodiment. The input links **230** and the output links **240** are shown as uni-directional. Preferably, the transmission format is a frame format, and includes frames **210** transmitted across the input/output link. In a preferred implementation of the first embodiment, the frames are transmitted back-to-back, in which case the frame boundaries can be identified by a frame receiver using a start-of-frame indicator (frame delimiter) **250**, and does not require end-of-frame indicators.

Each message frame defines a number of message transmission timeslots **260**, wherein each message transmission timeslot **260** is used to transmit a message including empty messages. The position of each message transmission timeslot **260** is fixed relative to the frame boundary. A receiver

6

does not need any transmission overhead per message to identify the message boundaries within received message frames.

The specific format of the different message types depends on the specific utilization of the N×N switching apparatus. In one particular implementation, the messages can be divided into three general categories or types. The first type of messages is packet/cell scheduling request command messages which has been generated by an input unit (e.g., a request to transfer a packet/cell just received at the input unit to a particular output unit), and forwarded to the message controller **220**, where they are processed and terminated by the packet/cell arbiter **350**. This message type is only transmitted on the message controller's input links.

The second type of messages is packet/cell acknowledge command messages which has been generated by the message controller's packet/cell arbiter **350**, and forwarded to an input unit for processing. This message type is only transmitted on the message controller's output links.

The third type of messages is messages which are generated by input/output units **200**, and are switched between input/output units by being transparently switched across the message controller **220**. This category also includes messages which are copied and replicated inside the message controller **220**, and then being transmitted out of the message controller in multiple copies on different output links **240**. This message type is transmitted on the message controller's input and output links.

The first embodiment can allocate the input link's message transmission bandwidth between the first and third message type by pre-assigning each of the message transmission timeslots **260** per input message frame **250** for one of the two message types. The optimal ratio between available input link transmission bandwidth for these two message types depends on the specific implementation of the switching apparatus, and can be modified as needed to suit that particular implementation.

The first embodiment allocates the output link's message transmission bandwidth between the second and third message type by pre-assigning each of the message transmission timeslots **260** per output message frame **250** for one of the two message types. The optimal ratio between available output link transmission bandwidth for these two message types depends on the specific implementation of the switching apparatus, and can be modified as needed to suit that particular implementation.

In one possible implementation of the first embodiment, the typical message transmission timeslot **260** size may vary depending on the type of message, whereby the message transmission according to the first embodiment can therefore be optimized by defining individual sizes message transmission timeslots for each of the corresponding message types.

In a preferred implementation of the first embodiment, the size of the message transmission timeslot **260** matches the corresponding message size in the switching apparatus embodiment, such that padding of the message information to match the size of the message to the message transmission timeslot size can be avoided.

FIG. **3** shows a block diagram of the message controller **130** as incorporated in the first embodiment of the present invention. Message frames arrive at the input links **330**, and incoming messages are forwarded by the message frame parser **320** to either the message switch **340** or the packet/cell arbiter **350**, depending on the individual message type. Messages are transmitted in message frames on the output links **300**, and the message frames are generated by the message

frame generator **310**. The message frame generator **310** receives frames from both the message switch **340** and the packet/cell arbiter **350**.

FIG. **4** is a flow diagram showing the functional operation of the message frame parser **320** shown in FIG. **3**. A message frame arrives at the input link in step **400**. In step **410**, the message frame parser **320** will one by one identify the next message transmission timeslot in the message frame. In step **420**, the method determines if the next message transmission timeslot is empty. If No, the flow goes to step **430**; if Yes, the flow goes to step **460**. In step **430**, it is determined whether or not the non-empty message transmission timeslot contains a message destined for the packet/cell arbiter **350**. If No, in step **440** the message is forwarded to the message switch **340**; if Yes, in step **450** the message is forwarded to the packet/cell arbiter **350**. In step **460**, a determination is made as to whether or not the last message transmission timeslot in the message has been processed; if Yes, the process returns to step **400** to wait for another arbiter frame to arrive at the input link, and if No, the process returns to step **410** to identify the next message transmission timeslot in the current arbiter frame.

By such a method of message processing, incoming messages destined for the packet/cell arbiter **350** are forwarded to the packet/cell arbiter **350**, and messages destined for an input/output unit **100/190** are forwarded to the message switch **340**. When all message transmission timeslots in an incoming message frame have been processed, the message frame parser **340** waits for the arrival of the next message frame.

FIG. **5** is a flow diagram showing the functional operation of the message frame generator **310** shown in FIG. **3**. In step **500**, new message frames begin transmitting, whereby the messages frames are generated one-by-one. In step **510**, a next message transmission timeslot in the message frame is identified. In step **520**, it is determined whether or not the message transmission timeslot identified in step **510** is allocated for a switched message. If Yes, in step **540** a message is inserted from the message switch into the outgoing arbiter frame; and if No, in step **530** a message from the packet/cell arbiter **350** is inserted into the outgoing arbiter frame.

In other words, the method determines whether the message transmission timeslot is pre-assigned for the packet/cell arbiter or for the message switch. When a message transmission timeslot is pre-assigned for the message switch, a message from the message switch is inserted into the outgoing message frame in step **540**. When a message transmission timeslot is pre-assigned for the message switch **340**, but a message is not available from the message switch, an empty message is inserted into the outgoing message frame.

When a message transmission timeslot is pre-assigned for the packet/cell arbiter **350**, a message from the packet/cell arbiter **350** is inserted into the outgoing message frame in step **530**. When a message transmission timeslot is pre-assigned for the packet/cell arbiter **350**, but a message is not available from the packet/cell arbiter **350**, an empty message is inserted into the outgoing message frame. Step **550** determines whether or not this is the last message transmission timeslot in the message frame; if Yes the process returns to step **500**, and if No the process goes to step **510** to identify the next message transmission timeslot in the message frame.

FIG. **6** shows a block diagram showing components of a message switch **340** shown in FIG. **3**. The message switch **340** is integrated into the message controller **130**, and provides message switching between the message controller's input links to the message controller's output links.

The message switch **340** includes a message scheduler **600**, a message crossbar **660**, one message input queue **650** per input link, and one message output queue **610** per output link.

The message scheduler **600** determines when messages are switched from input message queues **650** to output message queues **620** via the message crossbar **660**, and updates the message crossbar switching configuration accordingly every scheduling cycle.

The message crossbar **660** provides connectivity from any input message queue to any output message queue, and is capable of broadcasting from any input message queue **650** to all of the output message queues **620** (or to any particular subset thereof).

In a preferred implementation, the message scheduler **600** implements four parallel arbiters (not shown) per output message queue **620**:

One arbiter selects between even numbered input message queues **650** in fixed ascending order.

One arbiter selects between even numbered input message queues **650** in fixed descending order.

One arbiter selects between odd numbered input message queues **650** in fixed ascending order.

One arbiter selects between odd numbered input message queues **650** in fixed descending order.

Arbitration is preferably only performed on the input message queues head-of-line message, and each input message queue can forward one message into the message crossbar per scheduling cycle. The four arbiters implemented per output message queue are capable of switching one or two messages originating from even numbered links plus one or two messages originating from odd numbered links to each output message queue per scheduling cycle. When the incoming messages are evenly distributed between even and odd numbered input links, the message scheduler is capable of switching up to four messages to each output message queue per scheduling cycle.

A message can be switched as a unicast message or as a broadcast message. Broadcast switching is preferably performed spatially, meaning that the switching may be performed across multiple scheduling cycles. Once the message has been switched to all output message queues, it is removed from the input message queue head-of-line position. In a best case scenario, complete broadcast can be performed in a single scheduling cycle.

The methodology of different embodiments of the present invention has now been described above. The following will describe different options and approaches for implementing the invention.

Another embodiment of a switching apparatus incorporating a message switching method and apparatus is similar to the first embodiment shown in FIG. **1**, except that the output units **190** are not connected to the message controller **130**. In this embodiment, the input units (but not the output units) can switch messages between themselves via the message controller **130**. The input units can also forward packet/cell arbitration request command messages to the message controller's packet/cell arbiter **350**, and receive and process packet/cell arbitration acknowledge command messages generated by the message controller's packet/cell arbiter.

A switching apparatus incorporating a message switching method and apparatus of the foregoing embodiments includes a packet/cell switch element **160** for switching packets/cells between input and output units. The present invention can be incorporated with any packet/cell switch element that can provide switching of packets and/or cells between the input and output units. As one example, the packet/cell switch element can be implemented as a structure that includes a

single stage of parallel switch devices. Another possible implementation is a structure that includes multiple stages of switch devices.

The switch device for the packet/cell switch element can take on any of a number of forms that provide switching of packet and/or cells between the switch device input and output. Exemplary switch devices include crossbar switch devices, output buffered switch devices, crosspoint buffered switch devices, and switch devices embodied as described in U.S. patent application Ser. No. 10/898,540, entitled "Network interconnect Crosspoint Switching Architecture and Method", which is incorporated in its entirety herein by reference.

The message switch **340** can take on any number of forms that is able to provide switching of messages between the message controller's inputs and outputs. FIG. 6 depicts just one possible implementation of the message switch **340**.

Numerous variations of the message switch may exist. For example, the message switch may be implemented to have the number of arbiters per output to be less or more than the four arbiters per output in the message switch.

Another possible variation of the message switch implementation may correspond to having each arbiter select in a round robin fashion between inputs, instead of fixed ascending/descending order selection between inputs. The message switch can be implemented to have each arbiter select between all inputs instead of only between odd or even numbered inputs.

Yet another variation of the message switch implementation may correspond to having the input and/or output message queues implement multiple priority queuing levels, and/or where the message scheduler schedules message across the crossbar according to these priorities, instead of a single message priority.

Still another variation of the message switch implementation has an output buffered structure where each output buffer accepts simultaneously arriving messages from all inputs.

Further, while FIG. 2 depicts a preferred transmission format for transmitting message between input/output units and the message controller, other formats exist. One contemplated variation is to incorporate a transmission format where the message frames are not guaranteed to be transmitted back-to-back. In this variation, the boundary of the message frame is identified using both start-of-frame and end-of-frame identifiers per frame.

In another variation of the message transmission format, the message transmission timeslots in the message frames are not pre-assigned for specific message types, but are instead dynamically assigned by the message frame generator to the different message types. A field embedded in each message is used to identify the type of the message being transmitted in a message frame's message transmission timeslot.

Thus, apparatuses and methods have been described according to the present invention. Many modifications and variations may be made to the techniques and structures described and illustrated herein without departing from the spirit and scope of the invention. Accordingly, it should be understood that the methods and apparatus described herein are illustrative only and are not limiting upon the scope of the invention. Further, one or more aspects as described can be combined in any given system or method. Still further, one or more embodiments may be implemented in hardware, e.g., by a schematic design or a hardware description language (HDL), and/or implemented in a programmable logic device (FPGA/CPLD) or an ASIC, and/or they can be implemented in hardware using discrete hardware devices. Alternatively, one or more embodiments may be implemented in software.

The foregoing description is merely illustrative in nature and is in no way intended to limit the disclosure, its application, or uses. The broad teachings of the disclosure can be implemented in a variety of forms. Therefore, while this disclosure includes particular examples, the true scope of the disclosure should not be so limited since other modifications will become apparent upon a study of the drawings, the specification, and the following claims. As used herein, the phrase at least one of A, B, and C should be construed to mean a logical (A or B or C), using a non-exclusive logical OR. It should be understood that one or more steps within a method may be executed in different order (or concurrently) without altering the principles of the present disclosure.

What is claimed is:

1. A switching apparatus, comprising:

a switch element; and

a message controller separate from the switch element;

a plurality of input units configured to i) provide packets to the switch element, and ii) provide, to the message controller, messages corresponding to respective ones of the packets; and

a plurality of output units configured to i) receive, from the message controller, the messages, and ii) receive, from the switch element, the packets to be output from the switching apparatus, wherein

the switch element is configured to transfer each of the packets between the input units and the output units, and the message controller is configured to transfer each of the messages corresponding to the respective ones of the packets between the input units and the output units.

2. The switching apparatus of claim 1, wherein the message controller is configured to respectively transmit acknowledgement messages to the input units in response to receiving selected ones of the messages from the input units, wherein the acknowledgement messages direct the packets from the input units to the output units via the switch element.

3. The switching apparatus of claim 2, wherein the selected ones of the messages include scheduling request messages that request transfer of the packets between the input units and the output units.

4. The switching apparatus of claim 2, wherein the message controller includes an arbiter configured to generate the acknowledgement messages.

5. The switching apparatus of claim 1, wherein the message controller includes a message switch configured to transfer the messages between the input units and the output units.

6. The switching apparatus of claim 1, wherein the messages have a frame format.

7. The switching apparatus of claim 1, wherein the message controller is configured to i) assign message transmission timeslots to the messages, and ii) transmit the messages to the output units according to the message transmission timeslots.

8. The switching apparatus of claim 7, wherein the message controller is configured to insert, in the messages, indicators corresponding to a start of each of the messages and/or an end of each of the messages.

9. A method for operating a switching apparatus, the method comprising:

from a plurality of input units,

providing packets to a switch element, and

providing, to a message controller separate from the switch element, messages corresponding to respective ones of the packets;

using the switch element, transferring each of the packets between the input units and a plurality of output units; and

using the message controller, transferring each of the messages corresponding to the respective ones of the packets between the input units and the output units.

10. The method of claim **9**, further comprising respectively transmitting acknowledgement messages from the message controller to the input units in response to receiving selected ones of the messages from the input units, wherein the acknowledgement messages direct the packets from the input units to the output units via the switch element. 5

11. The method of claim **10**, wherein the selected ones of the messages include scheduling request messages that request transfer of the packets between the input units and the output units. 10

12. The method of claim **10**, wherein the acknowledgement messages are generated using an arbiter. 15

13. The method of claim **9**, further comprising using a message switch to transfer the messages between the input units and the output units.

14. The method of claim **9**, wherein the messages have a frame format. 20

15. The method of claim **9**, further comprising, using the message controller, i) assigning message transmission timeslots to the messages, and transmitting the messages to the output units according to the message transmission timeslots. 25

16. The method of claim **15**, further comprising, using the message controller, inserting indicators corresponding to a start of each of the messages and/or an end of each of the messages. 30

* * * * *